



# Genuine atomic multicast on Apache Kafka

E. Steinmacher, D. Pasadakis, E. De Lima Batista, O. Schenk, F. Pedone, & P. Eugster

Università della Svizzera Italiana, Institute of Computing, Lugano, Switzerland.  
Università della Svizzera Italiana, Computer Systems Institute, Lugano, Switzerland.



## Scaling throughput within Apache Kafka

- **Apache Kafka** is an event streaming platform that guarantees total order (TO) within a partition but not across partitions.<sup>[1]</sup>
- ➔ **Proof of concept:** Expand the Apache Kafka framework with TO multicast across partitions.



### Partitioning

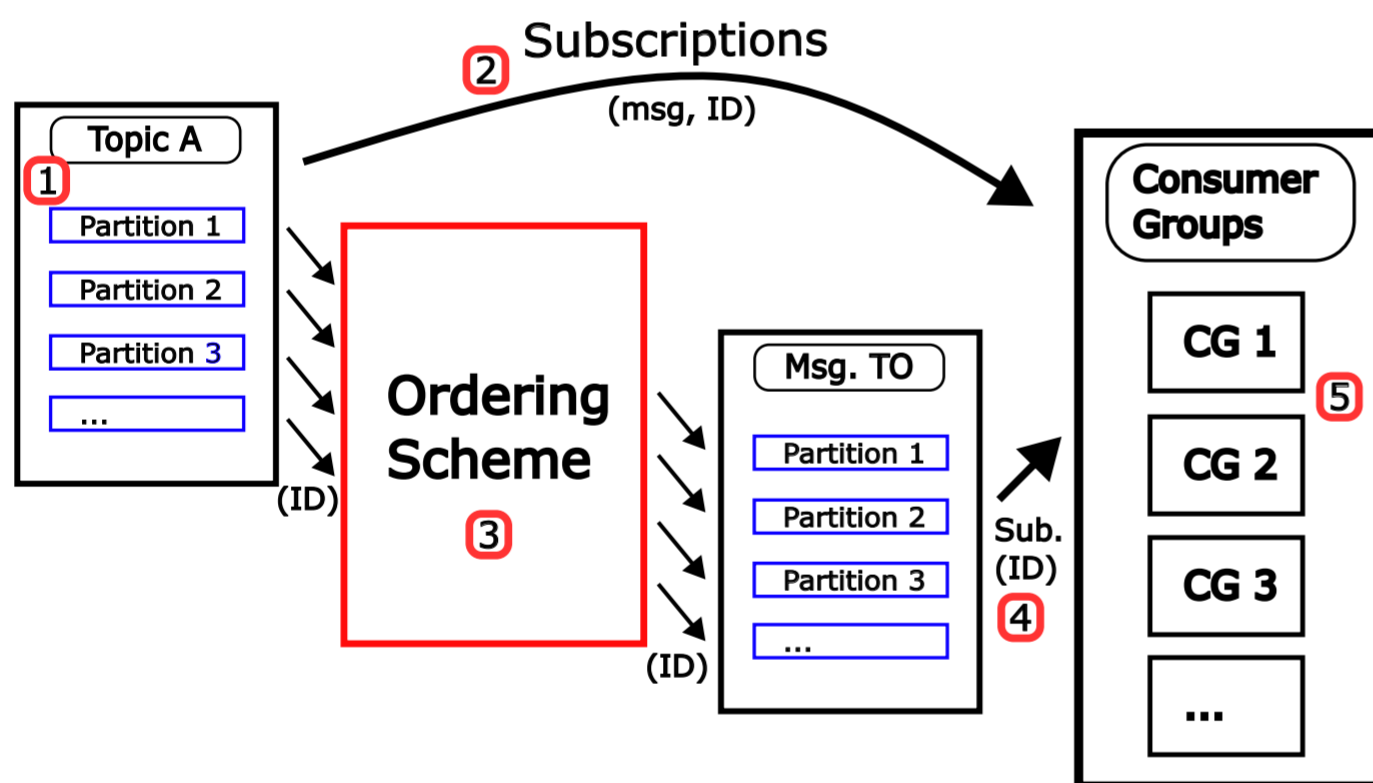
- Producers can create multiple data streams. Each data stream is input to a partition within a topic.
- Partitioning provides scalability. However, if multiple consumer groups are subscribed to the same subset of partitions, they do not consume the messages in the same order.
- ➔ Need of a new **ordering scheme**.

### Model assumptions

- Apache Kafka replicates partitions, thus providing fault-tolerance.
- Channels are FIFO and reliable.

### Setup of data pipeline:

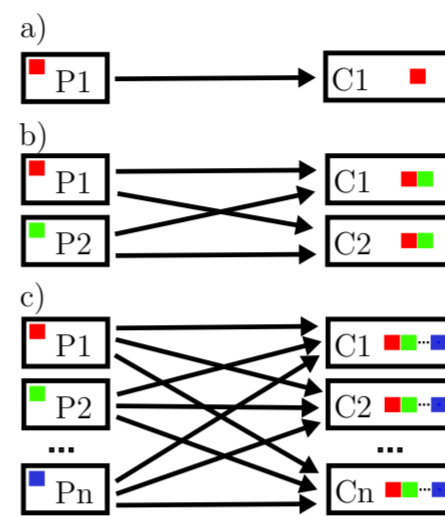
1. Message-ID tuples (msg, ID) are stored within the partitions of topic A.
2. Consumer groups are subscribed to topic A. Consumers within their group are assigned to partitions of topic A. ➔ Subscription link: Sends (msg, ID) tuples.
3. The hierarchical ordering scheme, orders the IDs and produces them to the topic "Msg. TO".
4. Consumer groups are also subscribed to "Msg. TO".
5. Consumers deliver a received message (msg, ID) as soon as they also received the ID via the ordering scheme.



## Numerical results

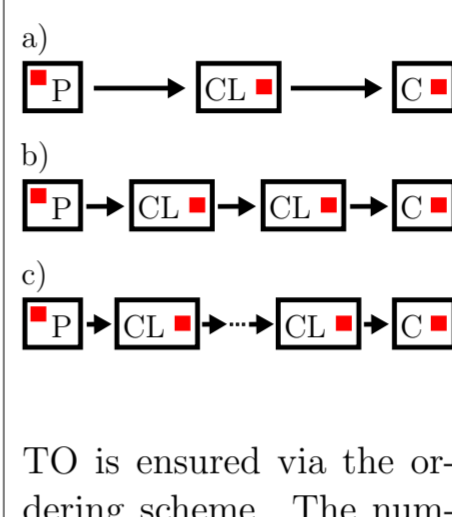
- Data sets are generated according to their number of messages  $\text{num\_msg} \in [10^2, 10^4]$  and message size  $\text{msg\_size} \in [10^4, 10^7]$  [B]. The ratio  $r = \log(\text{num\_msg})/\log(\text{msg\_size})$  is characteristic for a data set and is used to label.

### Kafka without TO



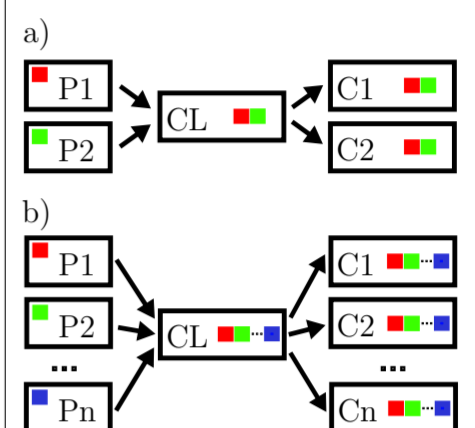
All consumers (C) are subscribed to all partitions (P). a) 1 P, 1 C; b) 2 P, 2 C; c) n P, n C, where  $n \in \{4, 8, 16, 32\}$ .

### TO w. tree height

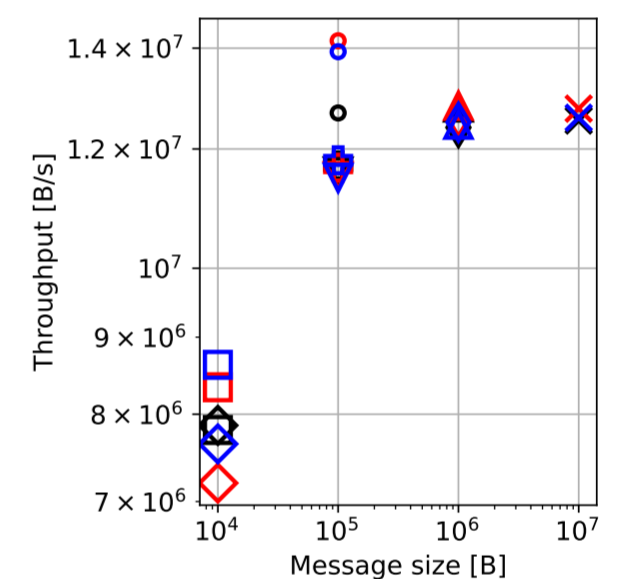
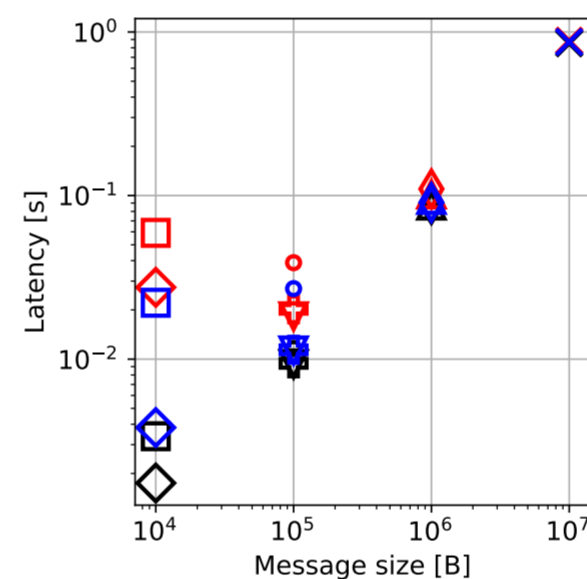


TO is ensured via the ordering scheme. The number of collectors in between is a) 1, b) 2, c)  $n \in \{4, 8, 16, 32\}$ .

### TO w. bottleneck size



TO is ensured via the ordering scheme. The size of the bottleneck ranges from 1, via a) 2, to b)  $n \in \{4, 8, 16, 32\}$ .

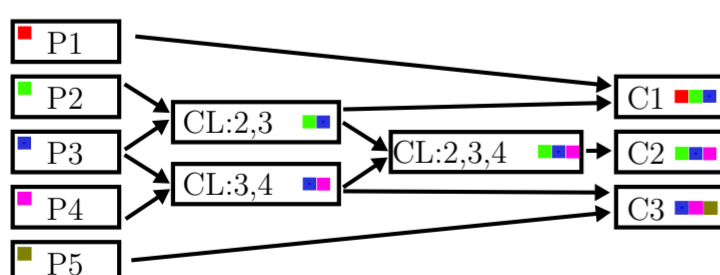


- Each data point represents the mean of all configurations (either 1, 2, ..., n P & C or 1, 2, ..., n CL) of a topology (Kafka without TO, TO w. tree height, TO w. bottleneck size), since different configurations show no clear trend.
- For small  $\text{msg\_size}$  the latency is increased by 5x to 10x, due to Kafka's batching for small and many messages. There is no latency increase for large  $\text{msg\_size}$ .
- There is no throughput decrease when the Kafka is expanded with TO across partitions.

So far this scheme is not suitable for latency critical applications when  $r \geq 10$ . Kafka's intrinsic properties have to be further explored, especially its batching of messages.

## Ordering Schemes & Implementation

### Hierarchical scheme



Toy example: Consumer (C) C1 consumes partitions (P) P1, P2 and P3. C2 consumes P2, P3 and P4. C3 consumes P3, P4, P5. The arrangement of the tree is an optimization problem. <sup>[2]</sup>

- Consumers are subscribed to partitions. All their common interests have to follow a total order.
- Msg-IDs are sent from partitions to consumers via layers formed by collectors (CL). The collectors resemble common interests.
- A consumer receives its order by having subscribed to the collectors that resemble its common interest ➔ ordering of the msg-IDs.

### Merge algorithm

- If the lower level collectors share one or multiple partitions, the streams have to be merged, such that the lower level streams order is preserved.

#### Algorithm 1 Merge algorithm

**Input:** Input streams  $s_i, i \in [1, \dots, N]$ , set of shared partitions  $S$   
**Output:** Merged stream preserving TO

```

1:  $S_{count} \leftarrow 0$ ;
2:  $s_{1,count} \leftarrow 0, \dots, s_{N,count} \leftarrow 0$ ;
3:  $q_{s_1} \leftarrow \emptyset, \dots, q_{s_N} \leftarrow \emptyset$ ;
4: for msg in consumer do
5:   Match msg to its  $q$  and  $s_{count}$ 
6:   if  $s_{count} == S_{count}$  and  $\text{msg.partition} \notin S$  then
7:      $\text{send}(\text{msg})$ 
8:   else
9:      $q.append(\text{msg})$ 
10:    if  $\text{msg.partition} \in S$  then
11:       $s_{count} \leftarrow s_{count} + 1$ 
12:    if  $s_{i,count} > S_{count} \forall i \in [1, 2, \dots, N]$  then
13:       $\forall q: \text{send}(\text{msg})$  until  $\text{msg.partition} \in S$ 
14:       $\text{send}(\text{msg})$ , where  $\text{msg.partition} \in S$  ▷ del. same msg
15:    in all other  $q$ 
16:     $S_{count} \leftarrow S_{count} + 1$ 
17:     $\forall q: \text{send}(\text{msg})$  until  $\text{msg.partition} \in S$  or  $q == \emptyset$ 

```

### Implementation

- Python interface using `kafka-python`.
- Work in progress:
  - Move to distributed system (e.g. Confluent Cloud).
  - Implement scheme using `kafka streams`.

### References

- [1] The Apache Software Foundation. <https://kafka.apache.org>, 2024.
- [2] Paulo Coelho, Tarcisio Ceolin Junior, Alysson Bessani, Fernando Dotti, and Fernando Pedone. Byzantine fault-tolerant atomic multicast. *48th Annual IEEE/IFIP International Conference on Dependable Systems and Networks*, 2018.

Student researcher project by Ekkehard Steinmacher.  
[ekkehard.steinmacher@usi.ch](mailto:ekkehard.steinmacher@usi.ch)